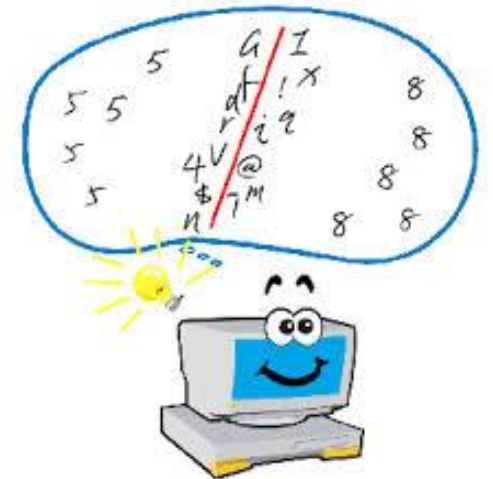# Quadratic Classification

Reporter : 王士銘

# Outline

- Classification

- Minimum Error – Maximum A Posterior

- Estimation of parameters
    - Maximum Likelihood
    - Bayesian Estimation

- Discernment Function

- Conclusion

# Classification

- Given an object, tell which class it belongs to
- Step -
    - Learn from training objects with class labels
    - Classifying testing objects without class labels
- Assumption - Each class of object follow a probability distribution with some parameters
- Machine learning theory guarantees small training error leads to small testing error
- The object of classifying algorithm is thus minimizing training error

# Notations

- $c$ – number of classes
- $D$ – training data
- $D_i$ - subset of training data with elements belonging to class $i$
- $N, n_i$ – size of $D$ and $D_i$
- $x, y$ – feature vector and class label representing a training object
- $\hat{y}$ - underground class label
- $\theta_i$ – parameters of the probability distribution of class $i$
- $\widehat{\theta_i}$ - maximum likelihood estimator for $\theta_i$

# Minimum Error (Testing)

- Consider 2 classes. Given training data D, to classify a testing object **x**

- $P(error|x)=\begin{cases} 1-P(y_1|x,D) \ , if \ we \ choose \ y=1 \\ 1-P(y_2|x,D) \ , if \ we \ choose \ y=2 \end{cases}$

- Hence, choose class 1 if $P(y_1|\pmb{x},D) > P(y_2|\pmb{x},D)$

- Minimizing error rate is equivalent to maximizing $P(y|\pmb{x},D)$, posterior probability

- $P(y_i|\pmb{x},D) = \dfrac{P(\pmb{x}|y_i,D\ )P(\pmb{y_i}|D)}{P(\pmb{x},D)} = \dfrac{P(\pmb{x}|y_i,D_i\ )P(\pmb{y_i})}{P(\pmb{x},D)}$

- Independence : $D_j, j \neq i$ gives no more information to $P(x|y_i,D_i)$

- Independence : $D, y_i$ are independent

- We need to estimate $P(\pmb{x}|y_i,D)$ and $P(y_i)$ for $i = 1,2 \dots, c$

# Estimation of Parameters(Training)

- need to estimate $P(\boldsymbol{x}|y_i, D)$ and $P(y_i)$ for $i = 1, 2 \ldots, c$

- $P(y_i) = \dfrac{n_i}{N}$

- $P(\boldsymbol{x}|y_i, D_i) = \int P(x|\theta_i, D_i)P(\theta_i|D_i)d\theta_i = \int P(x|\theta_i)P(\theta_i|D_i)d\theta$

- $P(\theta_i|D_i) = \dfrac{P(D_i|\theta_i)P(\theta_i)}{P(D_i)}$

- Maximum Likelihood Estimation(MLE) : estimate $\widehat{\theta_i} = \arg \underset{\theta}{max} \, \mathrm{P(D}|\theta_i)$
  and let $P(\theta_i|D_i) = \delta_{\widehat{\theta_i}}(\theta) \Rightarrow P(x|y_i, D_i) = P(x|\theta_i)$

- Bayesian estimation : $P(\theta_i|D_i) = \dfrac{P(D_i|\theta_i)P(\theta_i)}{P(D_i)} = \alpha\{\Pi_{k=1}^{n_i}P(x_k|\theta_i)\}P(\theta_i)$

- If both $P(\theta_i)$ and $P(x|\theta_i)$ are Gaussian, then $P(\theta_i|D_i)$ is Gaussian

# Maximum Likelihood Estimation(Training)

- $P(\theta_i|D_i) = \frac{P(D_i|\theta_i)P(\theta_i)}{P(D_i)} = \alpha\{\Pi_{k=1}^{n_i}P(x_k|\theta_i)\}P(\theta_i)$

- Estimate $\widehat{\theta_i} = \arg\max_{\theta} P(D|\theta_i)$ and let $P(\theta_i|D_i) = \delta_{\widehat{\theta_i}}(\theta)$

- Define log-likelihood function $l(\theta_i) \equiv lnP(D_i|\theta_i) = \Sigma_{k=1}^{n}\ln P(x_k|\theta_i)$

- $\widehat{\theta_i} = \arg\max_{\theta_i} l(\theta_i) \Rightarrow \nabla_{\theta_i}l(\widehat{\theta_i})=0$

  where $\nabla_{\boldsymbol{\theta}} = \begin{bmatrix} \frac{\partial}{\partial\theta_1} \\ \vdots \\ \frac{\partial}{\partial\theta_d} \end{bmatrix}$

# Maximum Likelihood Estimation(Training)

- Consider normal distribution for each class $P(\boldsymbol{x}|\theta_i) = N(\boldsymbol{x}|\mu_i, \Sigma_i)$

- $l(\theta_i) = \Sigma_{k=1}^{n_i} lnN(x_k|\mu_i, \Sigma_i) =$
  $\frac{n_i}{2} \ln|\Sigma| + \frac{-1}{2} \Sigma_{k=1}^{n_i} (x_k - \mu_i)^T \Sigma_i^{-1} (x_k - \mu_i) + constant$

- Taking the derivative over $\mu_i$ and set it to zero :

- $\Sigma_{k=1}^{n_i} \Sigma_i^{-1} (x_i - \mu_i) = 0$

- $\widehat{\mu_i} = \frac{1}{n_i} \Sigma_{k=1}^{n_i} x_i$

# Maximum Likelihood Estimation (Training)

- $l(\theta_i) = \frac{n_i}{2} \ln |\Sigma| + \frac{-1}{2} \Sigma_{k=1}^{n_i} (x_k - \mu_i)^T \Sigma_i^{-1} (x_k - \mu_i) + constant$

$$\propto \frac{n_i}{2} l \, n|\Sigma| + \frac{-1}{2} \Sigma_{k=1}^{n_i} Trace\big(\Sigma^{-1}(x_k - \mu_i)^T (x_k - \mu_i)\big)$$

$$= \frac{n_i}{2} l \, n|\Sigma| + \frac{-1}{2} Teace(\Sigma_{k=1}^{n_i} \Sigma^{-1}(x_k - \mu_i)^T (x_k - \mu_i))$$

- Taking the derivative over $\Sigma_i$ and set it to zero using

  1) $\frac{\partial}{\partial \Sigma} \ln |\Sigma| = \Sigma^{-T}$

  2) $\frac{\partial}{\partial \Sigma} Tr(\Sigma A) = \frac{\partial}{\partial \Sigma} Tr(A\Sigma) = A^T$

- $\widehat{\Sigma}_i = \frac{1}{n_i} \Sigma_{k=1}^{n_i} (x_k - \widehat{\mu}_i)(x_k - \widehat{\mu}_i)^T$

# Discriminant Function (Testing)

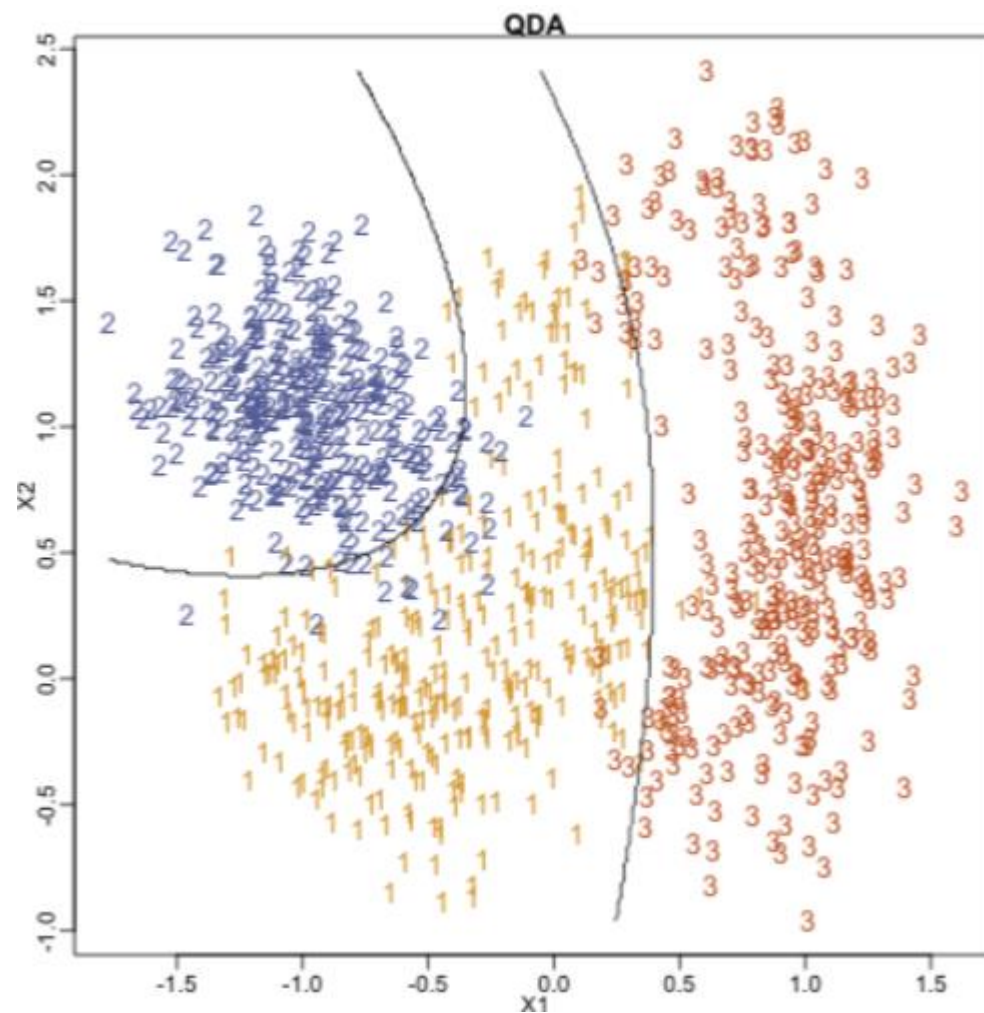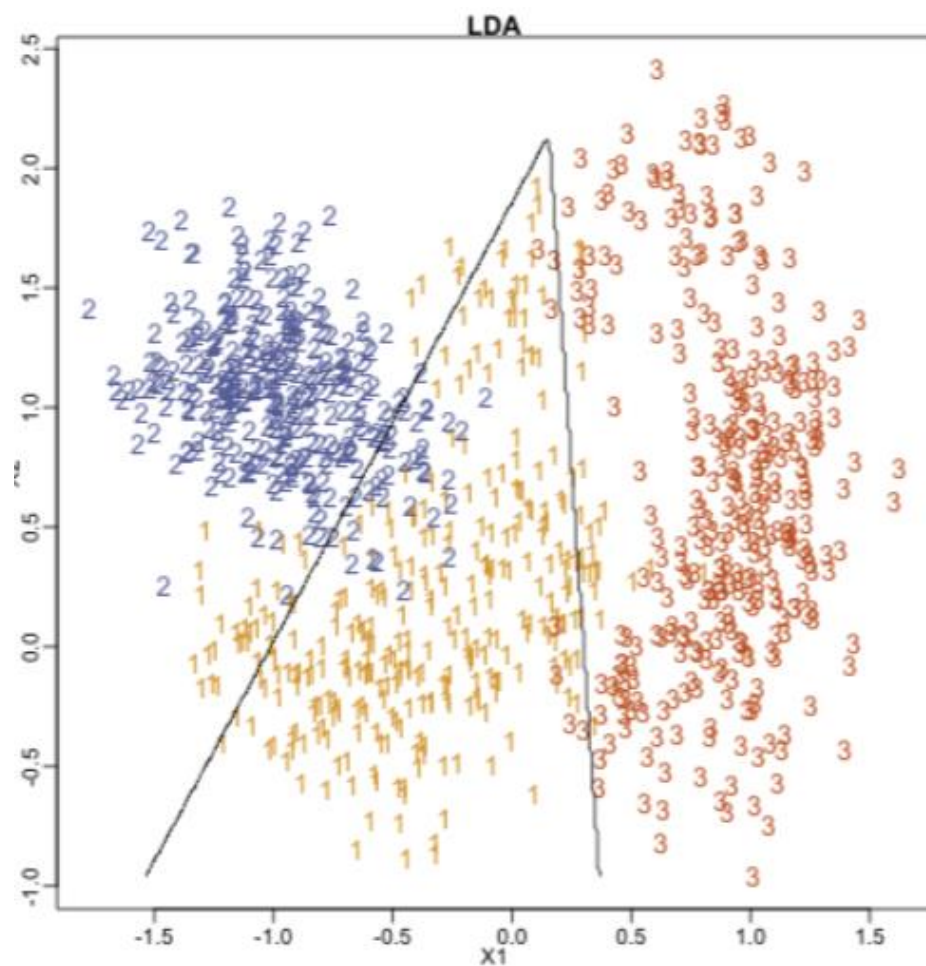- Suppose we model distribution of each class as multivariate Gaussian and estimate the mean $\mu_i$ and covariance matrix $\Sigma_i$ by MLE

- $P(y_i|\boldsymbol{x}) = \dfrac{P(\boldsymbol{x}|\widehat{\theta}_i)P(\boldsymbol{y}_i)}{P(\boldsymbol{x})} = \dfrac{N(\boldsymbol{x}|\widehat{\mu_i},\widehat{\Sigma_i})P(\boldsymbol{y}_i)}{P(\boldsymbol{x})}$

- Choose class $i$ if $P(y_i|\boldsymbol{x},D) > P(y_j|\boldsymbol{x},D)$ for $j \neq i$

- Define discrimination function $\delta_{ij}(\boldsymbol{x}) = \ln(P(y_i|\boldsymbol{x},D)/P(y_j|\boldsymbol{x},D))$

- Choose $i$ instead of $j$ if $\delta_{ij}(x) > 0$

# Discriminant Function (Testing)

- Discrimination function $\delta_{ij}(\boldsymbol{x}) = \ln(P(y_i|\boldsymbol{x}, D)/P(y_j|\boldsymbol{x}, D))$ determines the boundary for class $i$ and $j$

- If we assume $\widehat{\Sigma}_i = \Sigma$ for $i = 1, 2, \ldots, c$

- $\delta_{ij}(\boldsymbol{x}) = (\widehat{\mu_i} - \widehat{\mu_0})^T \widehat{\Sigma}^{-1}\boldsymbol{x} - \frac{1}{2}(\widehat{\mu_j} - \widehat{\mu_i})^T \widehat{\Sigma}^{-1}(\widehat{\mu_j} - \widehat{\mu_i}) + \ln\frac{P(y_i)}{P(y_j)}$

  $= \boldsymbol{w}^T\boldsymbol{x} + \boldsymbol{w_0}$ **->** linear discernment function

- If we assume different $\widehat{\Sigma}_i = \Sigma$ for $i = 1, 2, \ldots, c$

- $\delta_{ij}(\boldsymbol{x})$ will be quadratic function of $\boldsymbol{x}$ **->** quadratic discernment function

# Discriminant Function (Testing)

# Conclusion

- A traditional algorithm but has not bad performance
- Chinese handwritten recognition with Gabor filter feature

| | |
|---|---|
| QDA | 0.8724 |
| Neural Network | 0.8870 |
| SVM | 0.918945 |